

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2003-36146

(P2003-36146A)

(43) 公開日 平成15年2月7日(2003.2.7)

(51) Int.Cl. ⁷	識別記号	F I	データ* (参考)
G 0 6 F 3/06	3 0 5	G 0 6 F 3/06	3 0 5 C 5 B 0 1 8
	5 4 0		5 4 0 5 B 0 6 5
12/16	3 2 0	12/16	3 2 0 L

審査請求 未請求 請求項の数 4 O L (全 6 頁)

(21) 出願番号 特願2001-221566(P2001-221566)

(22) 出願日 平成13年7月23日(2001.7.23)

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田豊河台四丁目6番地

(72) 発明者 金子 誠司

神奈川県小田原市中里322番地2号 株式

会社日立製作所RAIDシステム事業部内

(74) 代理人 100093492

弁理士 鈴木 市郎 (外1名)

Fターム(参考) 5B018 GA02 GA04 HA14 MA12

5B065 BA01 CA11 CA30 CC08 EA03

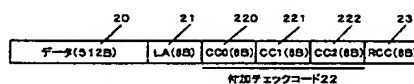
EA12 EA37 ZA15

(54) 【発明の名称】 ディスクアレイ制御方式

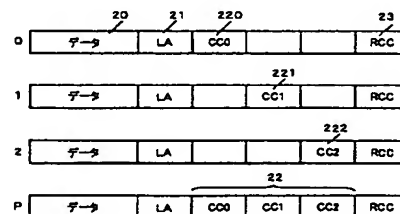
(57) 【要約】

【課題】 エラー検出能力の低いディスク制御部で、誤って旧データが読まれている障害を検出するとともに、付加チェックコードによってデータのエラーをも検出できるようにすること。

【解決手段】 複数のディスク装置を使って構成したディスクアレイ装置において、パリティグループを構成する複数のデータディスクに、各データディスクに書き込まれた各データに基づいて生成された各付加チェックコード220、221、又は222をそれぞれ書き込むとともに、パリティグループを構成するパリティディスクに、全ての各付加チェックコード220、221、及び222を書き込み、複数データディスクの1つと前記パリティディスクの付加チェックコードを読み出してエラーチェックを行うディスクアレイ制御方式。



(1)



(2)

図2 ディスクデータフォーマット

(2) 開2003-36146 (P2003-3chA)

【特許請求の範囲】

【請求項1】 複数のディスク装置を使って構成したディスクアレイ装置において、

パリティグループを構成する複数のデータディスクに、各データディスクに書き込まれた各データに基づいて生成された各付加チェックコードをそれぞれ書き込むとともに、

前記パリティグループを構成するパリティディスクに、全ての各付加チェックコードを書き込むことを特徴とするディスクアレイ制御方式。

【請求項2】 複数のディスク装置を使って構成したディスクアレイ装置において、

パリティグループを構成する複数のデータディスクに、各データディスクに書き込まれた各データに基づいて生成された各付加チェックコードをそれぞれ書き込むとともに、前記パリティグループを構成するパリティディスクに、全ての各付加チェックコードを書き込み、

前記複数データディスクの1つと前記パリティディスクの前記付加チェックコードを読み出してエラーチェックを行うことを特徴とするディスクアレイ制御方式。

【請求項3】 請求項2に記載のディスクアレイ制御方式において、

前記複数のデータディスクと前記パリティディスクに、書き込むべきデータのディスクアレイ内の位置情報をそれぞれ書き込み、誤った位置に書き込まれたデータのエラーチェックに前記位置情報を用いることを特徴とするディスクアレイ制御方式。

【請求項4】 請求項3に記載のディスクアレイ制御方式において、

パリティグループを構成するデータディスクとパリティディスクに、各データと各位置情報と各付加チェックコードとから生成されたチェックコードをディスクデータフォーマットの最後に書き込むことを特徴とするディスクアレイ制御方式。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、ディスクアレイ装置の障害検出方式に関し、特に、低価格で障害検出能力の低いディスク装置を組み合わせて、高信頼なディスクアレイ装置を構成する技術に関する。

【0002】

【従来の技術】一般的に、ディスクは磁気記録媒体や光記録媒体を使用しているために、円板上の一部分が読み書きできなくなる障害が発生する。これをセクタ障害という。セクタ障害の原因としては、円板上の傷や磁性体の劣化等が考えられる。ディスク装置では、ECC (Error Correcting Code) を付加することで、ある程度のセクタ障害に対しては、復旧することが可能である。また、付加するECC符号の種類により、セクタ中の何ビットまでのエラーを訂正すること

が可能であるかが変化する。

【0003】低価格のドライブでは、上記のセクタ障害が起きる確率が比較的高いため、複数のディスクを並列に動作させることで高速制御を実現し、また、パリティと呼ぶ冗長データをパリティディスクと呼ぶ特定のディスクに格納することにより、万一、データを格納する1台のディスクが故障しても、他のディスクとパリティディスクのパリティとから故障したディスクのデータを再現することができ、耐ディスク障害信頼性を高めることができるRAIDと呼ばれるディスク制御の方法が提案された。

【0004】このRAIDに関しては、「A Case for Redundant Arrays of Inexpensive Disks (RAID)」: In Proc. ACM SIGMOD, June 1988 (カリフォルニア大学バークレー校発行) に詳しい。RAIDは、そのパリティの格納の方法によりレベル1からレベル5があるが、現在広く耐障害性を持つディスクアレイで用いられているのはレベル1とレベル5である。

【0005】冗長構成を有する上記RAID5構成を取るディスクアレイ装置では、セクタ障害が発生した場合、同一パリティグループに属する残りのディスク装置によってデータを復元し、交代領域への書き込みを行うことにより障害データを回復させていた。

【0006】ここで、図3に示す4つのディスクへのデータ配置を例に取って、レイド5構成について述べる。と、ディスクアレイ制御装置に転送されてきたデータを0, 1, 2のように例えば512バイトの固定長に分けてディスク1, 2, 3に転送するとともに、パリティ部Pとしてデータ0, 1, 2のエクシクルーシブOAを取って512バイトの固定長を形成してディスク4に転送する。そして、このデータ0, 1, 2とパリティPが同一パリティグループを形成しているのである。図3の矢印に示すように、データ0, 1, 2, パリティP, データ0, 1, 2, パリティP, ...と続くディスク1~4へのデータ配置はレイド5の規格である。

【0007】そして、RAID5構成を取るひとつのディスクからデータを読み出そうとする場合、パリティグループを構成する全データディスクとパリティディスクからデータを読み出してエラーの検出をすることは可能であるが、複数のディスクからデータを読み出すため大きな性能低下を招く。

【0008】このため、このようなRAID装置の場合、所定のデータにエラー検出用のコードを付け、その両方をデータとして書き込み、読み出しをおこなうことにより単体ディスクの読み出しでエラー検出が行える様に構成するのが通常である。即ち、データ0, 1, 2には、前述した512バイトの生データに加えて、ディスクのロジカルアドレスLAとチェックビットCHを付加して

(3) 開2003-36146 (P2003-3ch(A))

いる。そして、このLAやCHによって単体ディスクでのエラー検出を行っている(図2を参照すると、図2のデータ20が生データに対応し、図2のLA21がロジカルアドレスに対応し、RCC23がチェックビットに対応するものである)。

【0009】また、特開平10-171608号公報には、データ部に対してアドレス情報と時間的要素を表す情報とを付加することが開示されている。このようなデータフォーマットを採用することで、磁気ディスク装置のキャッシュレジスタの古いデータを誤って読み出した場合でも磁気ディスク装置用データに変換する際にデータにタイムスタンプを追加記憶しておくので、データの読み出し又は書き込みの際にデータの異常を検出できるようになっている。

【0010】

【発明が解決しようとする課題】従来技術として説明した技術は、主にディスクの媒体に対する書き込み及び読み出し時のセクタ障害に代表されるエラーの検出、及びインターフェース関連のエラー検出とその回復を目的とする手段である。

【0011】一方、近年ディスクドライブはパーソナルコンピュータ用を中心に著しい価格の低下が進んだが、それにとめない低価格化を重視してメディア(記録媒体)の障害に比して発生頻度の低い内部制御系のチェックの省略が行われたため、ディスクの制御部でのエラー検出能力は相対的に低くなっている。この場合でも、ディスク内部で発生したエラーが、例えばデータ系のビット化けなどの単純なデータ誤りとして書き込まれるならば、従来技術で説明したRAIDデータのエラー検出符号によって検出可能であるが、制御部用のデータを保持するディスク内のバッファに障害が発生した場合、書き込みもうとしたデータを誤った位置に書き込む誤動作が発生し得る。

【0012】この種のエラーでは、書き込まれた先のデータ(誤った位置に書き込まれたデータ)も、上書きされずに残ったデータ(本来、上書きされて消去されるべき旧データ)も正しいチェックコードを持つ場合に、特に、書き込まれるべきデータが誤った場所に書き込まれたことによって書き込まれないで残った旧データは、これ自体単独では矛盾していないものなので、これが誤っていることを検出する手段が従来の技術では存在しない。

【0013】敷衍すると、512Bの生データ、位置情報のLA、チェックビットのCH、からなるデータ部と、生データから生成されたパリティデータ、LA及びCHからなるパリティ部と、で構成する従来のRAIDの場合、パリティ部(P)と全てのデータ部(O、1、2)との読み出しでエラーチェックすれば間違いなくエラー検出できるが、前述した読み出しの性能低下を考慮すると、データ部の読み出しでデータとチェックコード

CHのチェックで問題なければ正のデータとしてホストに上げるという取り扱いをする場合がある。このような場合において、データが誤った位置に書き込まれた際には、書き込まれるべき位置に存するデータの読み出し指令に対して(ディスクアレイ制御装置は当該位置に新しいデータが書き込まれていると認識しているから)、上書きで消去されずに残っていた旧データが読み出されることとなる。旧データはそれ自体ではデータとCHとが一致するものであるから、この旧データの読み出しにエラー検出が掛からない(LAも問題がないのであるから)こととなる。

【0014】本発明の目的は、互いに共通するデータに基づいた付加チェックコードをデータ部とパリティ部の両方に保持して、データ部とパリティ部の2つのディスクでその付加チェックコードの検出を行うことによって、誤って旧データが読まれている障害を検出するとともに、付加チェックコードによってデータのエラーをも検出できることにある。

【0015】

【課題を解決するための手段】本発明は、主として次のようなディスクアレイ制御方式を採用することによって、従来技術では為し得なかったエラー検出を実現する機能乃至作用を有するものである。

(1) パリティ部のパリティデータに加えて、パリティ付随情報として、RAIDを構成するデータ部に書き込んだ識別情報(図2に示すCC0、CC1又はCC2であり、データ20に基づいて生成されるもの)を付加チェックデータとして加えて書き込みを行う。

(2) 常にデータ部とパリティ部の二つのディスク単位でデータの読み出しを行い、前記データ部の識別情報とパリティ部のパリティ付随情報とを比較し、誤って古いデータが読まれている障害を検出する。

【0016】本発明においては、読み出しには必ず二個のディスクにアクセスするため、読みだし性能は多少低下するが、古いデータを誤って読み出す障害を検出することができる。

【0017】以上のような機能乃至作用を果たすために、本発明は主として次のような構成を採用する。複数のディスク装置を使って構成したディスクアレイ装置において、パリティグループを構成する複数のデータディスクに、各データディスクに書き込まれた各データに基づいて生成された各付加チェックコードをそれぞれ書き込むとともに、前記パリティグループを構成するパリティディスクに、全ての各付加チェックコードを書き込み、前記複数データディスクの1つと前記パリティディスクの前記付加チェックコードを読み出してエラーチェックを行うディスクアレイ制御方式。

【0018】

【発明の実施の形態】本発明の実施形態に係るディスクアレイ制御方式について、図1、図2及び図3を用いて

(4) 開2003-36146 (P2003-3ch+A)

以下詳細に説明する。図1は本発明の実施形態に係るディスクアレイ装置の概略的構成を示す図であり、図2は本実施形態におけるディスクに書き込むデータフォーマットを示す図であり、図3は本実施形態におけるディスク内のデータ配置を示す図である。

【0019】図1において、110はホストとなるサーバ10への通信路を示す。この通信路110には、一般にSCSIやファイバーチャネル等が使われる。図の12は、本実施形態における制御を実施するディスクアレイ制御装置であり、ホスト通信制御部121、装置制御を行うプロセッサ部122、データのキャッシングを行うキャッシュ部123、ディスクの制御を行うディスク制御部124、から構成される。また、ディスクアレイ制御装置12にはディスク13が接続される。

【0020】本実施形態では、図3を参照して、3個のデータに付きパリティデータを1つ用意し、4つのディスクでパリティグループを構成するRAID5構成を取るため、4の整数倍のディスク台数が必要であり、図3では簡単のため4台記載している。データとパリティ自体は、図3に示す様に分散されて配置されるため、データ用のディスクとパリティ用の専用ディスクに分かれてはいるわけではない。

【0021】図2は、RAID装置に書き込むデータのフォーマットを示す。データ部とパリティ部は同じフォーマットを使用する。図2の20はディスクの保持するデータ、図2の21はそのデータのディスク中の位置情報、図2の22は本発明の特徴である付加チェックコードデータ、図2の23はデータ、位置情報、付加チェックコードから生成されるチェックコードである。

【0022】本実施形態では、パリティディスクに書き込む付加チェックコードは、データ部の書き込みに用いる付加チェックコードと同一のもの（但し、データと位置情報から生成されたもの）を用い、またデータ部で用いる付加チェックコードは、パリティ部に用いる付加チェックコードを用いる。実際に書き込むデータに付随するチェックコードRCCは、データ、位置情報、付加チェックコードから生成された全体のデータに対するチェックコードである。ここで、付加チェックコードは前述したようにデータ20に基づいて生成されたチェックコードであることが特徴であり、CRC (Cyclic Redundancy Check) コードを用いても良い。いずれもデータ20に基づいたチェックコードであるから、この付加チェックコードを用いて、読み出しデータのエラーチェックも可能であり、更に、付加チェックコードの生成には既存のデータ及び/又は位置情報LAを利用して、特別な他の情報源を必要としたり、利用するものではない。また、付加チェックビットとしてデータを書き込んだ際に書き込みを一意に識別できる識別子を用いても良い。

【0023】図3は、本実施形態におけるディスク内の

データ配置を示す。パリティグループを構成する各ディスク30、31、32、33に、3つのデータ34、35、36及びパリティ37が図示の様に格納されるが、これは図示するように64K Bytes程度の細かい単位で別のディスクを使うようになっており、それによって負荷が高くなるパリティ部を全ディスクに分散させて負荷の均一化と性能向上を図っている。

【0024】本実施形態の実際の動作を以下詳細に記載する。本実施形態では、ホスト10側から通信路110経由でデータの書き込みが指定された場合、ホストから送られてきた書き込むべきデータを512 Bytes単位に分割して、まずキャッシュメモリ部123にバッファする。次に、制御プロセッサ122により、ディスク13に書くべきデータを図2のフォーマットに従って生成する。ここで、データ20は、ホストから送られてきたデータそのものであり、位置情報21は、内部で管理している書き込むべきデータのディスクアレイ内の位置である。これは8バイトからなり、上位2バイトは、ディスクアレイ制御装置12につながったディスクの通番（接続位置情報）であり、下位6 Bytesは、ディスク内のデータの書き込み開始セクタ番号を示す。

【0025】本実施形態では、書き込みの際には、対象とするデータが配置されるディスクと、パリティディスクの両方に対して同じ形式の図2に示す付加チェックコード22を持つデータを書き込む。通常のRAIDの場合にもチェックコードとしてパリティデータを書き込むので、本実施形態が従来フォーマットと異なる点は、前記付加チェックコード22を付随データとして書き込む点である。

【0026】この時、書き込むデータの付加チェックコードデータ22には、本実施形態のRAID構成に対応した、3つの付加チェックコードのフィールド220、221、222があるが、データ部の付加チェックコードは、自データの位置の付加チェックコードのみ更新、パリティデータの付加チェックコードは、書き込んだデータに対応する付加チェックコードのみ更新し、他の部分は元のままとする。

【0027】例えば、パリティグループの第一のデータ領域34にデータを書き込んだ場合（図3参照）、そのデータからパリティデータを再計算し、パリティデータの付加チェックコードをデータ領域34のフィールド220の位置に書き、データから生成したチェックコードをパリティデータ37のフィールド220の位置に書く。最後に、このデータ全体（データ20、LA21、付加チェックコード22）に対して所定の検出能力を持ったチェックコード23を生成し、データとして追加する。本実施形態では、8 Bytesのチェックコードを付けており、512 Bytesのデータにつき実際にディスクに書き込まれるのは552 Bytesである。この552 Bytesのデータをディスク制御装置12が

(5) 開2003-36146 (P2003-3ch-藤織)

ら実際のディスク13に書き込む。

【0028】以上説明したディスクデータフォーマットについて、図2の(2)を用いて再度説明する。ホストサーバ10から送られてきたデータは、ディスクアレイ制御装置12で、データ0、1、2のように固定長データに分割される。更に、パリティPはデータ0、1、2からエクシクルーシブオアを取って作成された512Bのデータである。これらデータ0、1、2、パリティPは、図3に示すようにディスク1〜4に配置される。ここで、データ0のデータフォーマットについて、本実施形態の特徴である付加チェックコードCC0が512Bのデータ、LA21に基づいて作成されてLA21に続いて書き込まれ、更に、チェックコードRCC23が512Bデータ、LA21及びCC0(220)に基づいて作成され書き込まれる。データ1及びデータ2についても、同様に付加チェックコードCC1(221)及びCC2(222)が図示のように書き込まれる。また、パリティPについては、データ20、LA21に続いて、CC0(220)、CC1(221)、CC2(222)が全て書き込まれる。

【0029】一方、本実施形態では読み出しの際の処理は通常と異なり、本実施形態に即して必ず対象とするデータが配置されるディスク30〜33の一つとパリティデータ37の両方を読み込む。読み出した付加チェックコード22が読み出したデータの対応する位置の付加チェックコードと一致しているかを確認してエラーの検出を行う。本実施形態の場合、エラー検出は以下の判定で行う。

【0030】(1) チェックコード23とデータ20が不一致の場合には、そのデータはエラーである。これはディスク内部の書き込み時のデータ系の障害によって正しいデータが書き込まれなかったものであると考えられる。

【0031】(2) チェックコード23とデータ20は一致しているが、読み出されたデータに埋め込まれた位置情報21と、そのデータ20が本来あるべき位置(ディスクアレイ制御装置が書き込みを指令した位置)が一致していなかった場合は、その不一致となっているデータはエラーである。これはディスク13の内蔵制御部障害によって誤った位置に書き込まれた先のが読み出されたケースである。

【0032】(3) 付加チェックコード22が一致していなかった場合には、まずパリティグループのデータを読み、そのデータから不一致の各データが誤っていたことを仮定した二通りのデータ再構築を行い、その結果と付加チェックコード22を比較して矛盾している側が誤ったデータである。これは、制御部障害によって書き込まれるべきデータが誤った位置に書き込まれたため、旧データが読み出されているケースである。即ち、図2の(2)の場合を例にすると、データ0とパリティPを読

み出して付加チェックコードが不一致の場合、データ0とパリティPのいずれかがエラーである。そこで、データ1とパリティP、データ2とパリティP、の組み合わせで付加チェックコードの一致をみて、データ0とパリティPのいずれかのエラーを検出する。データ0がエラーであると解ると、データ1、データ2及びパリティPからデータ0の回復を行う。

【0033】(4) 不一致やチェックコードとの不整合がなかった時にはデータは正しく、そのまま用いて良い。

【0034】上記判定によりデータのエラーが検出された場合には、改めてアレイを構成する全データを読み込み、そのデータから障害によって失われたデータを回復する。リカバリ時は、付加チェックコードデータを読み出したデータから再構築することを除けば通常のRAIDデータ回復手順と全く同じ処理である。

【0035】本実施形態では、データディスク3つに対してパリティディスク1つを置く構成を取ったが、チェックコード23のフィールド数を増減させることにより他のRAID構成に対応させることも容易に可能である。

【0036】本実施形態では、識別子としてCRC等の情報を用いたが、これはRAID装置としてのエラー回復処理が容易となるためである。例えば、データ部を書き込んだ書き込み通算番号をパリティ部に保持する様にしてもRAID5では、同様のエラー検出能力を得ることができる。この場合3台以下のパリティグループ構成では、どちらのデータが正しいかを知ることができないという難点があり、チェックコードを更に追加する必要がある。

【0037】

【発明の効果】本発明によれば、ディスク制御部でのエラー検出能力の低いディスクを使いながら、エラー検出能力を備え、かつ性能低下を抑えたRAID装置を構成することができる。

【図面の簡単な説明】

【図1】本発明の実施形態に係るディスクアレイ装置の概略的構成を示す図である。

【図2】本実施形態におけるディスクに書き込むデータフォーマットを示す図である。

【図3】本実施形態におけるディスク内のデータ配置を示す図である。

【符号の説明】

- 12 ディスクアレイ制御装置
- 13 ディスク
- 20 データ
- 21 位置情報
- 34、35、36 データ領域
- 37 パリティデータ領域
- 110 通信路

(6) 開2003-36146 (P2003-3ch8A)

121 ホスト通信制御部
122 制御プロセッサ

123 キャッシュメモリ部
124 ディスク制御部

【図1】

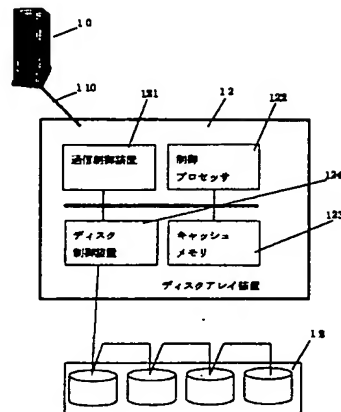


図1 ディスクアレイ構成

【図2】

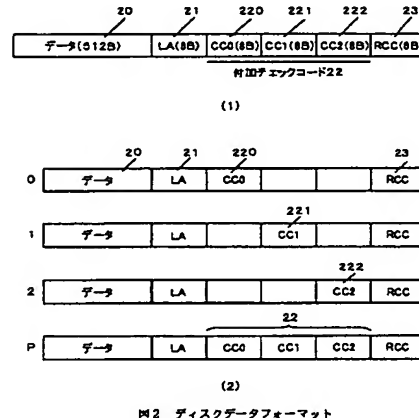


図2 ディスクデータフォーマット

【図3】

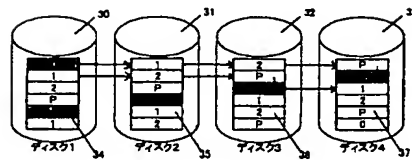


図3 ディスクデータ配座